

## Supplementary material for AAAI'19 Submission # 6931

### Appendix

**Proposition 3.** *The variance returned by Algorithm 1 is*

$$\sigma_{T,entire}^2(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}) - \sum_{i,j} \mathbf{z}_i^\top \left( \sum_l \mathbf{D}_l \mathbf{K}_{l,T} \mathbf{D}_l \right)^{-1} \mathbf{z}_j \quad (6)$$

where  $\mathbf{D}_j = \text{diag}([g_j(\mathbf{x}_1), \dots, g_j(\mathbf{x}_T)])$  and  $\mathbf{z}_i = \mathbf{D}_i \mathbf{k}_{j,T}(\mathbf{x}) g_j(\mathbf{x}) \in \mathbb{R}^T$ .

*Proof of Proposition 3.* According to Equation (4), the posterior covariance  $k_{T,entire}(\mathbf{x}, \mathbf{x}')$  in Algorithm 1 can be written as:

$$k_{T,entire}(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}') - \mathbf{k}_T(\mathbf{x})^\top \mathbf{K}_T^{-1} \mathbf{k}_T(\mathbf{x}') \quad (11)$$

By the decomposition assumption (Equation (1)), we have

$$k(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^J g_j(\mathbf{x}) k_j(\mathbf{x}, \mathbf{x}') g_j(\mathbf{x}'). \text{ Moreover,}$$

$$\begin{aligned} \mathbf{k}_T(\mathbf{x}) &= [k(\mathbf{x}_1, \mathbf{x}), \dots, k(\mathbf{x}_T, \mathbf{x})]^\top \\ &= \sum_{j=1}^J [g_j(\mathbf{x}_1) k_j(\mathbf{x}_1, \mathbf{x}) g_j(\mathbf{x}), \dots, g_j(\mathbf{x}_T) k_j(\mathbf{x}_T, \mathbf{x}) g_j(\mathbf{x})]^\top \\ &= \sum_{j=1}^J \mathbf{D}_j \mathbf{k}_{j,T}(\mathbf{x}) g_j(\mathbf{x}) \end{aligned} \quad (12)$$

where  $\mathbf{k}_{j,T}(\mathbf{x}) = [k_j(\mathbf{x}_1, \mathbf{x}), \dots, k_j(\mathbf{x}_T, \mathbf{x})]^\top$ .

The variance function  $\sigma_T^2(\mathbf{x})$  is just the value of covariance function with  $\mathbf{x}' = \mathbf{x}$ . Therefore, combining Equation (11) and (12), the variance can be written as:

$$\begin{aligned} \sigma_{T,entire}^2(\mathbf{x}) &= k_{T,entire}(\mathbf{x}, \mathbf{x}) \\ &= k(\mathbf{x}, \mathbf{x}) - \mathbf{k}_T(\mathbf{x})^\top \mathbf{K}_T^{-1} \mathbf{k}_T(\mathbf{x}) \\ &= k(\mathbf{x}, \mathbf{x}) - \sum_{i,j} g_i(\mathbf{x}) \mathbf{k}_{i,T}(\mathbf{x})^\top \mathbf{D}_i^\top \mathbf{K}_T^{-1} \mathbf{D}_j \mathbf{k}_{j,T}(\mathbf{x}) g_j(\mathbf{x}) \\ &= k(\mathbf{x}, \mathbf{x}) - \sum_{i,j} \mathbf{z}_i^\top \mathbf{K}_T^{-1} \mathbf{z}_j \end{aligned} \quad (13)$$

$$= k(\mathbf{x}, \mathbf{x}) - \sum_{i,j} \mathbf{z}_i^\top \left( \sum_l \mathbf{D}_l \mathbf{K}_{l,T} \mathbf{D}_l \right)^{-1} \mathbf{z}_j \quad (14)$$

with  $\mathbf{z}_i = \mathbf{D}_i \mathbf{k}_{j,T}(\mathbf{x}) g_j(\mathbf{x}) \in \mathbb{R}^n$  and equation (13) to (14) is coming from:

$$\mathbf{K}_T = [k(\mathbf{x}, \mathbf{x}')]_{\mathbf{x}, \mathbf{x}' \in A_T} + \text{diag}([\sigma^2(\mathbf{x}_t)]_{t \in [T]}) \quad (15)$$

$$\begin{aligned} &= \sum_{j=1}^J [g_j(\mathbf{x}) k_j(\mathbf{x}, \mathbf{x}') g_j(\mathbf{x}')]_{\mathbf{x}, \mathbf{x}' \in A_T} \\ &\quad + \text{diag}([g_j^2(\mathbf{x}_t) \sigma_j^2(\mathbf{x}_t)]_{t \in [T]}) \end{aligned} \quad (16)$$

$$\begin{aligned} &= \sum_j \mathbf{D}_j ([k_j(\mathbf{x}, \mathbf{x}')]_{\mathbf{x}, \mathbf{x}' \in A_T} + \text{diag}([\sigma_j^2(\mathbf{x}_t)]_{t \in [T]})) \mathbf{D}_j \\ &= \sum_j \mathbf{D}_j \mathbf{K}_{j,T} \mathbf{D}_j \end{aligned}$$

where the first kernel term from Equation (15) to (16) is derived by from definition. And in the latter term, from the decomposition assumption (1), the noise variance  $\sigma^2(\mathbf{x})$  of the target function  $f$  at point  $\mathbf{x}$  is the cumulative variance of the noise variance  $\sigma_j^2(\mathbf{x})$  of each individual function  $f_j$ , i.e.

$$\sigma^2(\mathbf{x}) = \sum_{j=1}^J g_j^2(\mathbf{x}) \sigma_j^2(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X}, j \in [J]$$

which explains the derivation from Equation 15 to 16.  $\square$

**Proposition 4.** *The variance returned by Algorithm 2 is*

$$\sigma_{T,decomp}^2(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}) - \sum_l \mathbf{z}_l^\top (\mathbf{D}_l \mathbf{K}_{l,T} \mathbf{D}_l)^{-1} \mathbf{z}_l \quad (7)$$

*Proof of Proposition 4.* In Algorithm 2, it runs GP regression to each function  $f_j(\mathbf{x})$  respectively. We can compute the corresponding posterior covariance function  $k_{j,T}$  by:

$$k_{j,T}(\mathbf{x}, \mathbf{x}') = k_j(\mathbf{x}, \mathbf{x}') - \mathbf{k}_{j,T}(\mathbf{x})^\top \mathbf{K}_{j,T}^{-1} \mathbf{k}_{j,T}(\mathbf{x}')$$

By Algorithm 2, the synthetic covariance of the target function  $f(\mathbf{x})$  is:

$$k_{T,decomp}(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^J g_j(\mathbf{x}) k_{j,T}(\mathbf{x}, \mathbf{x}') g_j(\mathbf{x}')$$

$$\begin{aligned} \sigma_{T,decomp}^2(\mathbf{x}) &= k_{T,decomp}(\mathbf{x}, \mathbf{x}) \\ &= \sum_{j=1}^J g_j(\mathbf{x}) k_{j,T}(\mathbf{x}, \mathbf{x}) g_j(\mathbf{x}) \\ &= \sum_{j=1}^J g_j(\mathbf{x}) k_j(\mathbf{x}, \mathbf{x}) g_j(\mathbf{x}) \\ &\quad - \sum_{j=1}^J g_j(\mathbf{x}) \mathbf{k}_{j,T}(\mathbf{x})^\top \mathbf{K}_{j,T}^{-1} \mathbf{k}_{j,T}(\mathbf{x}) g_j(\mathbf{x}) \\ &= k(\mathbf{x}, \mathbf{x}) - \sum_j \mathbf{z}_j^\top \mathbf{D}_j^{-1} \mathbf{K}_{j,T}^{-1} \mathbf{D}_j^{-1} \mathbf{z}_j \\ &= k(\mathbf{x}, \mathbf{x}) - \sum_j \mathbf{z}_j^\top (\mathbf{D}_j \mathbf{K}_{j,T} \mathbf{D}_j)^{-1} \mathbf{z}_j \end{aligned}$$

$\square$

### Proof of Theorem 6

**Theorem 6.** *The variance provided by decomposed Gaussian process regression (Algorithm 2) is less than or equal to the variance provided by Gaussian process regression (Algorithm 1), which implies the uncertainty by using decomposed Gaussian process regression is smaller.*

*Proof.* If we write  $B_l = D_l K_{l,T} D_l$ ,  $B$  is positive definite since it is the multiplication of positive definite matrix  $K_{l,T}$  and two  $D_l$  identical diagonal matrices. Using Proposition 3 and 4, the difference between Equation (6) and (7) can be written as:

$$\begin{aligned} & \sigma_{T,\text{entire}}^2(\mathbf{x}) - \sigma_{T,\text{decomp}}^2(\mathbf{x}) \\ &= \sum_l \mathbf{z}_l^\top (D_l K_{l,T} D_l)^{-1} \mathbf{z}_l - \sum_{i,j} \mathbf{z}_i^\top \left( \sum_l D_l K_{l,T} D_l \right)^{-1} \mathbf{z}_j \\ &= \sum_l \mathbf{z}_l^\top B_l^{-1} \mathbf{z}_l - \sum_{i,j} \mathbf{z}_i^\top \left( \sum_l B_l \right)^{-1} \mathbf{z}_j \\ &= \sum_l \mathbf{z}_l^\top B_l^{-1} \mathbf{z}_l - J \left( \frac{\sum_l \mathbf{z}_l}{J} \right)^\top \left( \frac{\sum_l B_l}{J} \right)^{-1} \left( \frac{\sum_l \mathbf{z}_l}{J} \right) \\ &= \sum_l h(B_l, \mathbf{z}_l) - J h(\bar{B}, \bar{\mathbf{z}}) \geq 0 \end{aligned}$$

where  $\bar{B} = \frac{\sum_l B_l}{J}$  and  $\bar{\mathbf{z}} = \frac{\sum_l \mathbf{z}_l}{J}$  are the average value. The last inequality comes from Jensen inequality and Lemma 5, which says the matrix-fractional function  $h$  is convex.  $\square$

### Proof of Theorem 7

In order to prove this, we follow the similar techniques as GPUCB (Srinivas et al. 2009), which is illustrated as follows:

**Lemma 9** (Modified version of Lemma 5.1 from Srinivas et al.). *Given  $f(\mathbf{x}) = \sum_{j=1}^J g_j(\mathbf{x}) f_j(\mathbf{x})$  (Definition 1), deterministic known functions  $g_j$  and unknown  $f_j \sim GP(0, k_j(\mathbf{x}, \mathbf{x}'))$ , pick  $\delta \in (0, 1)$  and set  $\beta_t = 2 \log(|\mathcal{X}| \pi_t / \delta)$ , where  $\sum_{t \geq 1} \pi_t^{-1} = 1, \pi_t > 0$ . Then, the  $\mu_{t-1}(\mathbf{x}), \sigma_{t-1}^2(\mathbf{x})$  returned by Algorithm 3 satisfy:*

$$|f(\mathbf{x}) - \mu_{t-1}(\mathbf{x})| \leq \beta_t^{1/2} \sigma_{t-1}(\mathbf{x}) \forall \mathbf{x} \in \mathcal{X}, t \geq 1$$

with probability  $1 - \delta$ .

*Proof.* Fix  $t \geq 1$  and  $\mathbf{x} \in \mathcal{X}$ . Conditioned on sampled points  $\{\mathbf{x}_1, \dots, \mathbf{x}_{t-1}\}$  and sampled values  $\{y_{1,j}, \dots, y_{t-1,j} \forall j \in [J]\}$ , the Bayesian property of decomposed GP regression (Algorithm 2) implies that the function value at point  $\mathbf{x}$  forms a Gaussian distribution with mean  $\mu_{t-1}(\mathbf{x})$  and variance  $\sigma_{t-1}^2(\mathbf{x})$ , i.e.  $f(\mathbf{x}) \sim N(\mu_{t-1}(\mathbf{x}), \sigma_{t-1}^2(\mathbf{x}))$ . Now, if  $r \sim N(0, 1)$ , then

$$\begin{aligned} Pr\{r > c\} &= e^{-c^2/2} (2\pi)^{-1/2} \int e^{-(r-c)^2/2 - c(r-c)} dr \\ &\leq e^{-c^2/2} Pr\{r > 0\} = (1/2) e^{-c^2/2} \end{aligned} \quad (17)$$

for  $c > 0$ , since  $e^{-c(r-c)} \leq 1$  for  $r \geq c$ . Therefore,  $Pr\{|f(\mathbf{x}) - \mu_{t-1}(\mathbf{x})| > \beta_t^{1/2} \sigma_{t-1}(\mathbf{x})\} \leq e^{-\beta_t^{1/2}}$ , using  $r = (f(\mathbf{x}) - \mu_{t-1}(\mathbf{x})) / \sigma_{t-1}(\mathbf{x})$  and  $c = \beta_t^{1/2}$ . Then apply the union bound to all  $\mathbf{x} \in \mathcal{X}$ :

$$|f(\mathbf{x}) - \mu_{t-1}(\mathbf{x})| \leq \beta_t^{1/2} \sigma_{t-1}(\mathbf{x}) \forall \mathbf{x} \in \mathcal{X}$$

holds with probability  $\geq 1 - |\mathcal{X}| e^{-\beta_t^{1/2}}$ . Choosing  $|\mathcal{X}| e^{-\beta_t^{1/2}} = \delta / \pi_t$  and using the union bound for  $t \in \mathbb{N}$ , the statement holds. For example, we can use  $\pi_t = \pi^2 t^2 / 6$ .  $\square$

The proof is almost the same as Theorem 5.1 in (Srinivas et al. 2009) except the Bayesian property of decomposed Gaussian process, where the Bayesian property of decomposed Gaussian process can be gotten from the Bayesian property of each individual function  $f_j$  and the linear combination of Gaussian distributions is still a Gaussian distribution, which implies the posterior belief after performing decomposed GP regression at a given point  $\mathbf{x}$  still form a Gaussian distribution with composed mean and variance.

**Lemma 10** (Modified version of Lemma 5.2 from Srinivas et al.). *Fix  $t \geq 1$ . If  $|f(\mathbf{x}) - \mu_{t-1}(\mathbf{x})| \leq \beta_t^{1/2} \sigma_{t-1}(\mathbf{x})$  for all  $\mathbf{x} \in \mathcal{X}$ , then the regret  $r_t$  is bounded by  $2\beta_t^{1/2} \sigma_{t-1}(\mathbf{x}_t)$ , where  $\mathbf{x}_t$  is the  $t$ -th choice of Algorithm 3.*

*Proof.* By definition of  $\mathbf{x}_t$ :  $\mu_{t-1}(\mathbf{x}_t) + \beta_t^{1/2} \sigma_{t-1}(\mathbf{x}_t) \geq \mu_{t-1}(\mathbf{x}^*) + \beta_t^{1/2} \sigma_{t-1}(\mathbf{x}^*) \geq f(\mathbf{x}^*)$ . Therefore,

$$\begin{aligned} r_t &= f(\mathbf{x}^*) - f(\mathbf{x}_t) \leq \beta_t^{1/2} \sigma_{t-1}(\mathbf{x}_t) + \mu_{t-1}(\mathbf{x}_t) - f(\mathbf{x}_t) \\ &\leq 2\beta_t^{1/2} \sigma_{t-1}(\mathbf{x}_t) \end{aligned}$$

$\square$

**Lemma 11** (Modified version of Lemma 5.3 from Srinivas et al.). *The information gain for the points selected can be expressed in terms of the predictive variances. If  $\mathbf{f}_{j,T} = \{f_j(\mathbf{x}_t)\}_{t \in [T]} \in \mathbb{R}^T$  and  $\mathbf{y}_{j,T} = \{y_{j,t}\}_{t \in [T]} \in \mathbb{R}^T$ :*

$$I(\mathbf{y}_{j,T} : \mathbf{f}_{j,T}) = \frac{1}{2} \sum_{t=1}^T \log(1 + \sigma^{-2} \sigma_{j,t-1}^2(\mathbf{x}_t))$$

where  $f_j(\mathbf{x}_t), y_{j,t}, \sigma_{j,t-1}^2$  follow the definition and derivation in Algorithm 3.

*Proof.* Directly follow by replacing all the  $f, y, \sigma$  by  $f_j, y_j, \sigma_j$  in the proof of Theorem 5.3 from Srinivas et al.  $\square$

**Theorem 7.** *Let  $\delta \in (0, 1)$  and  $\beta_t = 2 \log(|\mathcal{X}| t^2 \pi^2 / 6\delta)$ . Running decomposed GP-UCB (Algorithm 3) for a composed sample  $f(\mathbf{x}) = \sum_{j=1}^J g_j(\mathbf{x}) f_j(\mathbf{x})$  with bounded variance  $k_j(\mathbf{x}, \mathbf{x}) \leq 1$  and each  $f_j \sim GP(0, k_j(\mathbf{x}, \mathbf{x}'))$ , we obtain a regret bound of  $\mathcal{O}(\sqrt{T \log |\mathcal{X}| \sum_{j=1}^J B_j^2 \gamma_{j,T}})$  with high probability, where  $B_j = \max_{\mathbf{x} \in \mathcal{X}} |g_j(\mathbf{x})|$ . Precisely,*

$$Pr\{R_T \leq \sqrt{C_1 T \beta_T \sum_{j=1}^J B_j^2 \gamma_{j,T}} \forall T \geq 1\} \geq 1 - \delta \quad (8)$$

where  $C_1 = 8 / \log(1 + \sigma^{-2})$  with noise variance  $\sigma^2$ .

*Proof.* According to Lemma 11, we can take advantage of the individual information gain of each  $f_j(\mathbf{x})$ , which is

$$\begin{aligned} I_j(\mathbf{y}_{j,T}; f_j) &= \frac{1}{2} \sum_{t=1}^T \log(1 + \sigma^{-2} \sigma_{j,t-1}^2(\mathbf{x}_t)) \\ \gamma_{j,T} &= \max I_j(\mathbf{y}_{j,T}; f_j) \end{aligned}$$

Besides, we can also bound the total regret by the individual information gains as following:

$$\begin{aligned}
\sum_{j=1}^J B_j^2 I_j(y_{j,T}; f_j) &= \frac{1}{2} \sum_{j=1}^J B_j^2 \sum_{t=1}^T \log(1 + \sigma^{-2} \sigma_{j,t-1}^2(\mathbf{x}_t)) \\
&\geq \frac{1}{2} \sum_{j=1}^J B_j^2 \sum_{t=1}^T C_2^{-1} \sigma^{-2} \sigma_{j,t-1}^2(\mathbf{x}_t) \\
&\geq \frac{1}{2} C_2^{-1} \sigma^{-2} \sum_{j=1}^J \sum_{t=1}^T g_j^2(\mathbf{x}_t) \sigma_{j,t-1}^2(\mathbf{x}_t) \\
&\geq \frac{1}{2} C_2^{-1} \sigma^{-2} \sum_{t=1}^T \frac{r_t^2}{4\beta_t} \\
&\geq \frac{C_2^{-1} \sigma^{-2}}{8\beta_T} \sum_{t=1}^T r_t^2
\end{aligned}$$

where  $C_2 = \sigma^{-2} / \log(1 + \sigma^{-2}) \geq 1$ ,  $s^2 \leq C_2 \log(1 + s^2)$  for  $s \in [0, \sigma^{-2}]$  and  $\sigma^{-2} \sigma_{j,t-1}^2(\mathbf{x}_t) \leq \sigma^{-2} k_j(\mathbf{x}_t, \mathbf{x}_t) \leq \sigma^{-2}$ . Let  $C_1 = 8\sigma^2 C_2 = 8 / \log(1 + \sigma^{-2})$ . Applying Cauchy inequality gives us:

$$C_1 \beta_T T \sum_{j=1}^J B_j^2 I_j(y_{j,T}; f_j) \geq \left( \sum_{t=1}^T r_t \right)^2 = R_T^2$$

which implies a similar upper bound

$$R_T \leq \sqrt{C_1 T \beta_T \sum_{j=1}^J B_j^2 \gamma_{j,T}}$$

□

### Proof of Theorem 8

All the proofs in Theorem 7 apply except Lemma 9. Since the decomposition here is non-linear, therefore the composition of outcomes of Gaussian processes is no longer an outcome of Gaussian process, which prohibits us to have a nice Gaussian process property: function value  $f(\mathbf{x})$  does not form a Gaussian distribution. Due to the non-linearity, the distribution gets distorted, losing its original form with Gaussian distribution. Fortunately, if the partial derivatives of function  $g: \mathbb{R}^J \rightarrow \mathbb{R}$  (Definition 2) are bounded, then we can still perform a similar estimation and bound the distribution by a larger Gaussian distribution, which enables us to have a similar result.

**Lemma 12** (General Version with Definition 2 and Algorithm 4). *Given  $f(\mathbf{x}) = g(f_1(\mathbf{x}), \dots, f_J(\mathbf{x}))$  (Definition 2), deterministic known functions  $g$  and unknown  $f_j \sim GP(0, k_j(\mathbf{x}, \mathbf{x}'))$ , pick  $\delta \in (0, 1)$  and set  $\beta_t = 2 \log(|\mathcal{X}| J \pi_t / \delta)$ , where  $\sum_{t \geq 1} \pi_t^{-1} = 1, \pi_t > 0$ . Further assume the function  $g$  has bounded partial derivatives  $B_j = \max_{\mathbf{x} \in \mathcal{X}} |\nabla_j g(\mathbf{x})| \forall j \in [J]$ . Then, the  $\mu_{t-1}(\mathbf{x}), \sigma_{t-1}(\mathbf{x})$  returned by Algorithm 4 satisfy:*

$$|f(\mathbf{x}) - \mu_{t-1}(\mathbf{x})| \leq \beta_t^{1/2} \sigma_{t-1}(\mathbf{x}) \forall \mathbf{x} \in \mathcal{X}, t \geq 1$$

with probability  $1 - \delta$ .

*Proof.* The main problem here is the posterior distribution of  $f(\mathbf{x})$  is not a Gaussian distribution. But fortunately, the posterior distribution of each  $f_j(\mathbf{x})$  is still a Gaussian distribution with mean  $\mu_{j,t-1}(\mathbf{x})$  and variance  $\sigma_{j,t-1}^2(\mathbf{x})$  for any given  $\mathbf{x} \in \mathcal{X}$ . Then,

$$\begin{aligned}
&|f(\mathbf{x}) - \mu_{t-1}(\mathbf{x})| \\
&= |g(f_1(\mathbf{x}), \dots, f_J(\mathbf{x})) - g(\mu_{1,t-1}(\mathbf{x}), \dots, \mu_{J,t-1}(\mathbf{x}))| \\
&\leq \sum_{j=1}^J B_j |f_j(\mathbf{x}) - \mu_{j,t-1}(\mathbf{x})| \tag{18}
\end{aligned}$$

Applying the same argument in Lemma 9 to function  $f_j$ :

$$Pr\{|f_j(\mathbf{x}) - \mu_{j,t-1}(\mathbf{x})| > \beta_t^{1/2} \sigma_{j,t-1}(\mathbf{x})\} \leq e^{-\beta_t^{1/2}}$$

Then applying the union bound on  $j \in [J]$ , we get

$$\begin{aligned}
&|f(\mathbf{x}) - \mu_{t-1}(\mathbf{x})| \\
&\leq \sum_{j=1}^J B_j |f_j(\mathbf{x}) - \mu_{j,t-1}(\mathbf{x})| \\
&\leq \sum_{j=1}^J B_j \beta_t^{1/2} \sigma_{j,t-1}(\mathbf{x}) \\
&\leq \beta_t^{1/2} \sqrt{J \left( \sum_{j=1}^J B_j^2 \sigma_{j,t-1}^2(\mathbf{x}) \right)} \\
&= \beta_t^{1/2} \sigma_{t-1}(\mathbf{x})
\end{aligned}$$

with probability  $1 - J e^{-\beta_t^{1/2}}$ , where the last inequality is from Cauchy's inequality. Then apply union bound again to all  $\mathbf{x} \in \mathcal{X}$ , the above inequality yields:

$$|f(\mathbf{x}) - \mu_{t-1}(\mathbf{x})| \leq \beta_t^{1/2} \sigma_{t-1}(\mathbf{x}) \forall \mathbf{x} \in \mathcal{X}$$

with probability  $1 - |\mathcal{X}| J e^{-\beta_t^{1/2}}$ . Choosing  $|\mathcal{X}| J e^{-\beta_t^{1/2}} = \delta / \pi_t$  and using the union bound for  $t \in \mathbb{N}$ , the statement holds, i.e.  $\beta_t = 2 \log(|\mathcal{X}| J \pi_t / \delta)$ . Specifically, if we choose  $\pi_t = \pi^2 t^2 / 6$ , then it implies  $\beta_t = 2 \log(|\mathcal{X}| J t^2 \pi^2 / 6\delta)$ . □

**Theorem 8.** *By running generalized decomposed GP-UCB with  $\beta_t = 2 \log(|\mathcal{X}| J t^2 \pi^2 / 6\delta)$  for a composed sample  $f(\mathbf{x}) = g(f_1(\mathbf{x}), \dots, f_J(\mathbf{x}))$  of GPs with bounded variance  $k_j(\mathbf{x}, \mathbf{x}) \leq 1$  and each  $f_j \sim GP(0, k_j(\mathbf{x}, \mathbf{x}'))$ . we obtain a regret bound of  $\mathcal{O}(\sqrt{T \log |\mathcal{X}| \sum_{j=1}^J B_j^2 \gamma_{j,T}})$  with high probability, where  $B_j = \max_{\mathbf{x} \in \mathcal{X}} |\nabla_j g(\mathbf{x})|$ . Precisely,*

$$Pr\{R_T \leq \sqrt{C_1 T \beta_T \sum_{j=1}^J B_j^2 \gamma_{j,T}} \forall T \geq 1\} \geq 1 - \delta \tag{10}$$

where  $C_1 = 8 / \log(1 + \sigma^{-2})$  with noise variance  $\sigma^2$ .

*Proof.* Directly follow by the same proofs of Theorem 7 with Lemma 10, Lemma 11, and Lemma 12. □

**Remark 1.** In inequality (18), if we write  $Z_j = |f_j(\mathbf{x}) - \mu_{j,t-1}(\mathbf{x})|$ , where  $f_j(\mathbf{x}) - \mu_{j,t-1}(\mathbf{x})$  is sampled from a normal distribution with 0 mean and  $\sigma_{j,t-1}^2(\mathbf{x})$  (due to Gaussian process property). Then this  $Z_j$  is a random variable drawn from a half-normal distribution with parameter  $\sigma_j(\mathbf{x})$  (no longer the variance here).

*The summation of half-normal distributions can still be computed and bounded by a similar inequality like inequality (17). This can provide a constant ratio of improvement to the  $\beta_t$  exploration parameter, thus the regret bound as well. However it does not change the order of regret and sample complexity. Therefore we are not going to cover this here.*